

Sufficient Component Analysis for Supervised Dimension Reduction

Makoto Yamada
Gang Niu
Jun Takagi
Masashi Sugiyama

Tokyo Institute of Technology, Tokyo 152-8552, Japan

YAMADA@SG.CS.TITECH.AC.JP
GANG@SG.CS.TITECH.AC.JP
TAKAGI@SG.CS.TITECH.AC.JP
SUGI@CS.TITECH.AC.JP

Abstract

The purpose of sufficient dimension reduction (SDR) is to find the low-dimensional subspace of input features that is sufficient for predicting output values. In this paper, we propose a novel *distribution-free* SDR method called *sufficient component analysis* (SCA), which is computationally more efficient than existing methods. In our method, a solution is computed by iteratively performing dependence estimation and maximization: Dependence estimation is analytically carried out by recently-proposed *least-squares mutual information* (LSMI), and dependence maximization is also analytically carried out by utilizing the *Epanechnikov kernel*. Through large-scale experiments on real-world image classification and audio tagging problems, the proposed method is shown to compare favorably with existing dimension reduction approaches.

1. Introduction

The goal of *sufficient dimension reduction* (SDR) is to learn a transformation matrix \mathbf{W} from input feature \mathbf{x} to its low-dimensional representation \mathbf{z} ($= \mathbf{W}\mathbf{x}$) which has ‘sufficient’ information for predicting output value \mathbf{y} . SDR can be formulated as the problem of finding \mathbf{z} such that \mathbf{x} and \mathbf{y} are conditionally independent given \mathbf{z} (Cook, 1998; Fukumizu et al., 2009).

Earlier SDR methods developed in statistics community, such as *sliced inverse regression* (Li, 1991), *principal Hessian direction* (Li, 1992), and *sliced average variance estimation* (Cook, 2000), rely on the elliptical assumption (e.g., Gaussian) of the data, which may not be fulfilled in practice.

To overcome the limitations of these approaches, the *kernel dimension reduction* (KDR) was proposed (Fukumizu et al., 2009). KDR employs a kernel-based dependence measure, which does not require the elliptical assumption (i.e., distribution-free), and the solution \mathbf{W} is computed by a gradient method. Although KDR is a highly flexible SDR method, its critical weakness is the kernel function choice—the performance of KDR depends on the choice of kernel functions and the regularization parameter, but there is no systematic model selection method available. Furthermore, KDR scales poorly to massive datasets since the gradient-based optimization is computationally demanding. Another important limitation of KDR in practice is that there is no good way to set an initial solution—many random restarts may be needed for finding a good local optima, which makes the entire procedure even slower and the performance of dimension reduction unstable.

To overcome the limitations of KDR, a novel SDR method called *least-squares dimension reduction* (LSDR) was proposed recently (Suzuki & Sugiyama, 2010). LSDR adopts a squared-loss variant of mutual information as a dependency measure, which is efficiently estimated by *least-squares mutual information* (LSMI) (Suzuki et al., 2009). A notable advantage of LSDR over KDR is that kernel functions and its tuning parameters such as the kernel width and the regularization parameter can be naturally optimized by cross-validation. However, LSDR still relies on a computationally expensive gradient method and there is no good initialization scheme.

In this paper, we propose a novel SDR method called *sufficient component analysis* (SCA), which can overcome the computational inefficiency of LSDR. In SCA, the solution \mathbf{W} in each iteration is obtained *analytically* by just solving an eigenvalue problem, which significantly contributes to improving the computational efficiency. Moreover, based on the above analytic-form

solution, we develop a method to design a good initial value for optimization, which further reduces the computational cost and help obtain a good local optimum solution.

Through large-scale experiments using the *PASCAL Visual Object Classes (VOC) 2010* dataset (Everingham et al., 2010) and the *Freesound* dataset (The Freesound Project, 2011), we demonstrate the usefulness of the proposed method.

2. Sufficient Dimension Reduction with Squared-Loss Mutual Information

In this section, we formulate the problem of *sufficient dimension reduction* (SDR) based on *squared-loss mutual information* (SMI).

2.1. Problem Formulation

Let $\mathcal{X}(\subset \mathbb{R}^d)$ be the domain of input feature \mathbf{x} and \mathcal{Y} be the domain of output data¹ \mathbf{y} . Suppose we are given n independent and identically distributed (i.i.d.) paired samples,

$$D^n = \{(\mathbf{x}_i, \mathbf{y}_i) \mid \mathbf{x}_i \in \mathcal{X}, \mathbf{y}_i \in \mathcal{Y}, i = 1, \dots, n\},$$

drawn from a joint distribution with density $p_{\mathbf{xy}}(\mathbf{x}, \mathbf{y})$.

The goal of SDR is to find a low-dimensional representation \mathbf{z} ($\in \mathbb{R}^m$, $m \leq d$) of input \mathbf{x} that is sufficient to describe output \mathbf{y} . More precisely, we find \mathbf{z} such that

$$\mathbf{y} \perp\!\!\!\perp \mathbf{x} \mid \mathbf{z}, \quad (1)$$

meaning that, given the projected feature \mathbf{z} , the feature \mathbf{x} is conditionally independent of output \mathbf{y} .

In this paper, we focus on linear dimension reduction scenarios:

$$\mathbf{z} = \mathbf{W}\mathbf{x},$$

where \mathbf{W} is a transformation matrix. \mathbf{W} belongs to the *Stiefel manifold* $\mathbb{S}_m^d(\mathbb{R})$:

$$\mathbb{S}_m^d(\mathbb{R}) := \{\mathbf{W} \in \mathbb{R}^{m \times d} \mid \mathbf{W}\mathbf{W}^\top = \mathbf{I}_m\},$$

where $^\top$ denotes the transpose and \mathbf{I}_m is the m -dimensional identity matrix. Below, we assume that the reduced dimension m is known.

¹ \mathcal{Y} could be either continuous (i.e., regression) or categorical (i.e., classification). Multi-dimensional outputs (e.g., multi-task regression and multi-label classification) and structured outputs (such as sequences, trees, and graphs) can also be handled in the proposed framework.

2.2. Dependence Estimation-Maximization Framework

Suzuki & Sugiyama (2010) showed that the optimal transformation matrix that leads to Eq.(1) can be characterized as

$$\mathbf{W}^* = \operatorname{argmax}_{\mathbf{W} \in \mathbb{R}^{m \times d}} \text{SMI}(Z, Y) \text{ s.t. } \mathbf{W}\mathbf{W}^\top = \mathbf{I}_m. \quad (2)$$

In the above, $\text{SMI}(Z, Y)$ is the *squared-loss mutual information*:

$$\text{SMI}(Z, Y) := \frac{1}{2} \mathbb{E}_{p_{\mathbf{z}}, p_{\mathbf{y}}} \left[\left(\frac{p_{\mathbf{zy}}(\mathbf{z}, \mathbf{y})}{p_{\mathbf{y}}(\mathbf{y})p_{\mathbf{z}}(\mathbf{z})} - 1 \right)^2 \right],$$

where $\mathbb{E}_{p_{\mathbf{z}}, p_{\mathbf{y}}}$ denotes the expectation over the marginals $p_{\mathbf{z}}(\mathbf{z})$ and $p_{\mathbf{y}}(\mathbf{y})$. Note that SMI is the *Pearson divergence* from $p_{\mathbf{zy}}(\mathbf{z}, \mathbf{y})$ to $p_{\mathbf{z}}(\mathbf{z})p_{\mathbf{y}}(\mathbf{y})$, while the ordinary mutual information is the Kullback-Leibler divergence from $p_{\mathbf{zy}}(\mathbf{z}, \mathbf{y})$ to $p_{\mathbf{z}}(\mathbf{z})p_{\mathbf{y}}(\mathbf{y})$. The Pearson divergence and the Kullback-Leibler divergence both belong to the class of f -divergences, which shares similar theoretical properties. For example, SMI is non-negative and is zero if and only if Z and Y are statistically independent, as ordinary mutual information.

Based on Eq.(2), we develop the following iterative algorithm for learning \mathbf{W} :

- (i) **Initialization:** Initialize the transformation matrix \mathbf{W} (see Section 3.3).
- (ii) **Dependence estimation:** For current \mathbf{W} , an SMI estimator $\widehat{\text{SMI}}$ is obtained (see Section 3.1).
- (iii) **Dependence maximization:** Given an SMI estimator $\widehat{\text{SMI}}$, its maximizer with respect to \mathbf{W} is obtained (see Section 3.2).
- (iv) **Convergence check:** The above (ii) and (iii) are repeated until \mathbf{W} fulfills some convergence criterion².

3. Proposed Method: Sufficient Component Analysis

In this section, we describe our proposed method called the *sufficient component analysis* (SCA).

3.1. Dependence Estimation

In SCA, we utilize a non-parametric SMI estimator called *least-squares mutual information* (LSMI)

² In experiments, we used the criterion that the improvement of $\widehat{\text{SMI}}$ is less than 10^{-6} .

(Suzuki et al., 2009), which was shown to achieve the optimal convergence rate (Suzuki & Sugiyama, 2010). Here, we review LSMI.

3.1.1. BASIC IDEA

A key idea of LSMI is to directly estimate the *density ratio*,

$$w(\mathbf{z}, \mathbf{y}) = \frac{p_{zy}(\mathbf{z}, \mathbf{y})}{p_z(\mathbf{z})p_y(\mathbf{y})},$$

without going through density estimation of $p_{zy}(\mathbf{z}, \mathbf{y})$, $p_z(\mathbf{z})$, and $p_y(\mathbf{y})$. Here, the density ratio function $w(\mathbf{z}, \mathbf{y})$ is directly modeled by

$$w_\alpha(\mathbf{z}, \mathbf{y}) = \sum_{\ell=1}^n \alpha_\ell K(\mathbf{z}, \mathbf{z}_\ell) L(\mathbf{y}, \mathbf{y}_\ell), \quad (3)$$

where $K(\mathbf{z}, \mathbf{z}')$ and $L(\mathbf{y}, \mathbf{y}')$ are kernel functions for \mathbf{z} and \mathbf{y} , respectively.

Then, the parameter $\alpha = (\alpha_1, \dots, \alpha_n)^\top$ is learned so that the following squared error is minimized:

$$J_0(\alpha) = \frac{1}{2} \mathbb{E}_{p_z, p_y} [(w_\alpha(\mathbf{z}, \mathbf{y}) - w(\mathbf{z}, \mathbf{y}))^2].$$

J_0 can be expressed as

$$J_0(\alpha) = J(\alpha) + \text{SMI}(Z, Y) + \frac{1}{2},$$

where

$$J(\alpha) = \frac{1}{2} \alpha^\top \mathbf{H} \alpha - \mathbf{h}^\top \alpha,$$

$$\begin{aligned} H_{\ell, \ell'} &= \mathbb{E}_{p_z, p_y} [K(\mathbf{z}, \mathbf{z}_\ell) L(\mathbf{y}, \mathbf{y}_\ell) K(\mathbf{z}, \mathbf{z}_{\ell'}) L(\mathbf{y}, \mathbf{y}_{\ell'})], \\ h_\ell &= \mathbb{E}_{p_{zy}} [K(\mathbf{z}, \mathbf{z}_\ell) L(\mathbf{y}, \mathbf{y}_\ell)], \end{aligned}$$

and $\text{SMI}(Z, Y)$ is constant with respect to α . Thus, minimizing J_0 is equivalent to minimizing J .

3.1.2. COMPUTING THE SOLUTION

Approximating the expectations in \mathbf{H} and \mathbf{h} included in J by empirical averages, we arrive at the following optimization problem:

$$\min_{\alpha} \left[\frac{1}{2} \alpha^\top \widehat{\mathbf{H}} \alpha - \widehat{\mathbf{h}}^\top \alpha + \lambda \alpha^\top \mathbf{R} \alpha \right],$$

where a regularization term $\lambda \alpha^\top \mathbf{R} \alpha$ is included for avoiding overfitting, $\lambda (\geq 0)$ is a regularization parameter, \mathbf{R} is a regularization matrix, and, for $\mathbf{z}_i = \mathbf{W} \mathbf{x}_i$,

$$\widehat{H}_{\ell, \ell'} = \frac{1}{n^2} \sum_{i, j=1}^n K(\mathbf{z}_i, \mathbf{z}_\ell) L(\mathbf{y}_i, \mathbf{y}_\ell) K(\mathbf{z}_j, \mathbf{z}_{\ell'}) L(\mathbf{y}_j, \mathbf{y}_{\ell'}),$$

$$\widehat{h}_\ell = \frac{1}{n} \sum_{i=1}^n K(\mathbf{z}_i, \mathbf{z}_\ell) L(\mathbf{y}_i, \mathbf{y}_\ell).$$

Differentiating the above objective function with respect to α and equating it to zero, we can obtain an analytic-form solution:

$$\widehat{\alpha} = (\widehat{\mathbf{H}} + \lambda \mathbf{R})^{-1} \widehat{\mathbf{h}}. \quad (4)$$

Based on the fact that $\text{SMI}(Z, Y)$ is expressed as

$$\text{SMI}(Z, Y) = \frac{1}{2} \mathbb{E}_{p_{zy}} [w(\mathbf{z}, \mathbf{y})] - \frac{1}{2},$$

the following SMI estimator can be obtained:

$$\widehat{\text{SMI}} = \frac{1}{2} \widehat{\mathbf{h}}^\top \widehat{\alpha} - \frac{1}{2}. \quad (5)$$

3.1.3. MODEL SELECTION

Hyper-parameters included in the kernel functions and the regularization parameter can be optimized by cross-validation with respect to J .

More specifically, the samples $\mathcal{Z} = \{(\mathbf{z}_i, \mathbf{y}_i)\}_{i=1}^n$ are divided into K disjoint subsets $\{\mathcal{Z}_k\}_{k=1}^K$ of (approximately) the same size. Then, an estimator $\widehat{\alpha}_{\mathcal{Z}_k}$ is obtained using $\mathcal{Z} \setminus \mathcal{Z}_k$ (i.e., all samples without \mathcal{Z}_k), and the approximation error for the hold-out samples \mathcal{Z}_k is computed as

$$J_{\mathcal{Z}_k}^{(K\text{-CV})} = \frac{1}{2} \widehat{\alpha}_{\mathcal{Z}_k}^\top \widehat{\mathbf{H}}_{\mathcal{Z}_k} \widehat{\alpha}_{\mathcal{Z}_k} - \widehat{\mathbf{h}}_{\mathcal{Z}_k}^\top \widehat{\alpha}_{\mathcal{Z}_k},$$

where, for $|\mathcal{Z}_k|$ being the number of samples in the subset \mathcal{Z}_k ,

$$\begin{aligned} [\widehat{\mathbf{H}}_{\mathcal{Z}_k}]_{\ell, \ell'} &= \frac{1}{|\mathcal{Z}_k|^2} \sum_{(\mathbf{z}, \mathbf{y}), (\mathbf{z}', \mathbf{y}') \in \mathcal{Z}_k} K(\mathbf{z}, \mathbf{z}_\ell) L(\mathbf{y}, \mathbf{y}_\ell) \\ &\quad \times K(\mathbf{z}', \mathbf{z}_{\ell'}) L(\mathbf{y}', \mathbf{y}_{\ell'}), \\ [\widehat{\mathbf{h}}_{\mathcal{Z}_k}]_\ell &= \frac{1}{|\mathcal{Z}_k|} \sum_{(\mathbf{z}, \mathbf{y}) \in \mathcal{Z}_k} K(\mathbf{z}, \mathbf{z}_\ell) L(\mathbf{y}, \mathbf{y}_\ell). \end{aligned}$$

This procedure is repeated for $k = 1, \dots, K$, and its average $J^{(K\text{-CV})}$ is outputted as

$$J^{(K\text{-CV})} = \frac{1}{K} \sum_{k=1}^K J_{\mathcal{Z}_k}^{(K\text{-CV})}.$$

We compute $J^{(K\text{-CV})}$ for all model candidates, and choose the model that minimizes $J^{(K\text{-CV})}$.

3.2. Dependence Maximization

Given an SMI estimator $\widehat{\text{SMI}}$ (5), we next show how $\widehat{\text{SMI}}$ can be efficiently maximized with respect to \mathbf{W} :

$$\max_{\mathbf{W} \in \mathbb{R}^{m \times d}} \widehat{\text{SMI}} \quad \text{s.t. } \mathbf{W} \mathbf{W}^\top = \mathbf{I}_m.$$

We propose to use a truncated negative quadratic function called the *Epanechnikov kernel* (Epanechnikov, 1969) as a kernel for \mathbf{z} :

$$K(\mathbf{z}, \mathbf{z}_\ell) = \max \left(0, 1 - \frac{\|\mathbf{z} - \mathbf{z}_\ell\|^2}{2\sigma_z^2} \right).$$

Let $I(c)$ be the indicator function, i.e., $I(c) = 1$ if c is true and zero otherwise. Then, for the above kernel, $\widehat{\text{SMI}}$ can be expressed as

$$\widehat{\text{SMI}} = \frac{1}{2} \text{tr}(\mathbf{W} \mathbf{D} \mathbf{W}^\top) - \frac{1}{2},$$

where $\text{tr}(\mathbf{A})$ is the trace of matrix \mathbf{A} , and

$$\begin{aligned} \mathbf{D} = & \frac{1}{n} \sum_{i=1}^n \sum_{\ell=1}^n \hat{\alpha}_\ell(\mathbf{W}) I \left(\frac{\|\mathbf{W} \mathbf{x}_i - \mathbf{W} \mathbf{x}_\ell\|^2}{2\sigma_z^2} < 1 \right) \\ & \times L(\mathbf{y}_i, \mathbf{y}_\ell) \left[\frac{1}{m} \mathbf{I}_d - \frac{1}{2\sigma_x^2} (\mathbf{x}_i - \mathbf{x}_\ell)(\mathbf{x}_i - \mathbf{x}_\ell)^\top \right]. \end{aligned}$$

Here, by $\hat{\alpha}_\ell(\mathbf{W})$, we explicitly indicated the fact that $\hat{\alpha}_\ell$ depends on \mathbf{W} .

Let \mathbf{D}' be \mathbf{D} with \mathbf{W} replaced by \mathbf{W}' , where \mathbf{W}' is a transformation matrix obtained in the previous iteration. Thus, \mathbf{D}' no longer depends on \mathbf{W} . Here we replace \mathbf{D} in $\widehat{\text{SMI}}$ by \mathbf{D}' , which gives the following simplified SMI estimate:

$$\frac{1}{2} \text{tr}(\mathbf{W} \mathbf{D}' \mathbf{W}^\top) - \frac{1}{2}. \quad (6)$$

A maximizer of Eq.(6) can be analytically obtained by $(\mathbf{w}_1 | \cdots | \mathbf{w}_m)^\top$, where $\{\mathbf{w}_i\}_{i=1}^m$ are the m principal components of \mathbf{D}' .

3.3. Initialization of \mathbf{W}

In the dependence estimation-maximization framework described in Section 2.2, initialization of the transformation matrix \mathbf{W} is important. Here we propose to initialize it based on dependence maximization without dimensionality reduction.

More specifically, we determine the initial transformation matrix as $(\mathbf{w}_1^{(0)} | \cdots | \mathbf{w}_m^{(0)})^\top$, where $\{\mathbf{w}_i^{(0)}\}_{i=1}^m$ are

the m principal components of $\mathbf{D}^{(0)}$:

$$\begin{aligned} \mathbf{D}^{(0)} = & \frac{1}{n} \sum_{i=1}^n \sum_{\ell=1}^n \hat{\alpha}_\ell^{(0)} I \left(\frac{\|\mathbf{x}_i - \mathbf{x}_\ell\|^2}{2\sigma_x^2} < 1 \right) L(\mathbf{y}_i, \mathbf{y}_\ell) \\ & \times \left[\frac{1}{m} \mathbf{I}_d - \frac{1}{2\sigma_x^2} (\mathbf{x}_i - \mathbf{x}_\ell)(\mathbf{x}_i - \mathbf{x}_\ell)^\top \right], \end{aligned}$$

$$\hat{\boldsymbol{\alpha}}^{(0)} = (\widehat{\mathbf{H}}^{(0)} + \lambda \mathbf{R})^{-1} \hat{\mathbf{h}}^{(0)},$$

$$\begin{aligned} \widehat{\mathbf{H}}_{\ell, \ell'}^{(0)} = & \frac{1}{n^2} \sum_{i,j=1}^n K'(\mathbf{x}_i, \mathbf{x}_\ell) L(\mathbf{y}_i, \mathbf{y}_\ell) \\ & \times K'(\mathbf{x}_j, \mathbf{x}_{\ell'}) L(\mathbf{y}_j, \mathbf{y}_{\ell'}), \end{aligned}$$

$$\hat{\mathbf{h}}_\ell^{(0)} = \frac{1}{n} \sum_{i=1}^n K'(\mathbf{x}_i, \mathbf{x}_\ell) L(\mathbf{y}_i, \mathbf{y}_\ell),$$

$$K'(\mathbf{x}, \mathbf{x}_\ell) = \max \left(0, 1 - \frac{\|\mathbf{x} - \mathbf{x}_\ell\|^2}{2\sigma_x^2} \right).$$

σ_x is the kernel width and is chosen by cross-validation (see Section 3.1.3).

4. Relation to Existing Methods

Here, we review existing SDR methods and discuss the relation to the proposed SCA method.

4.1. Kernel Dimension Reduction

Kernel Dimension Reduction (KDR) (Fukumizu et al., 2009) tries to directly maximize the conditional independence of \mathbf{x} and \mathbf{y} given \mathbf{z} under a kernel-based independence measure.

The KDR learning criterion is given by

$$\begin{aligned} \mathbf{W}^* = & \underset{\mathbf{W} \in \mathbb{R}^{m \times d}}{\text{argmax}} \text{tr} \left[\tilde{\mathbf{L}} (\tilde{\mathbf{K}} + n\epsilon \mathbf{I}_n)^{-1} \right] \\ \text{s.t. } & \mathbf{W} \mathbf{W}^\top = \mathbf{I}_m, \end{aligned} \quad (7)$$

where $\tilde{\mathbf{L}} = \mathbf{\Gamma} \mathbf{L} \mathbf{\Gamma}$, $\mathbf{\Gamma} = \mathbf{I} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top$, $L_{i,j} = L(\mathbf{y}_i, \mathbf{y}_j)$, $\tilde{\mathbf{K}} = \mathbf{\Gamma} \mathbf{K} \mathbf{\Gamma}$, $K_{i,j} = K(\mathbf{z}_i, \mathbf{z}_j)$, and ϵ is a regularization parameter.

Solving the above optimization problem is cumbersome since the objective function is non-convex. In the original KDR paper (Fukumizu et al., 2009), a gradient method is employed for finding a local optimal solution. However, the gradient-based optimization is computationally demanding due to its slow convergence and it requires many restarts for finding a good local optima. Thus, KDR scales poorly to massive datasets.

Another critical weakness of KDR is the kernel function choice. The performance of KDR depends on the choice of kernel functions and the regularization parameter, but there is no systematic model selection

method for KDR available. Using the Gaussian kernel with its width set to the median distance between samples is a standard heuristic in practice, but this does not always work very well.

Furthermore, KDR lacks a good way to set an initial solution in the gradient procedure. Then, in practice, we need to run the algorithm many times with random initial points for finding good local optima. However, this makes the entire procedure even slower and the performance of dimension reduction unstable.

The proposed SCA method can successfully overcome the above weaknesses of KDR—SCA is equipped with cross-validation for model selection (Section 3.1.3), its solution can be computed analytically (see Section 3.2), and a systematic initialization scheme is available (see Section 3.3).

4.2. Least-Squares Dimensionality Reduction

Least-squares dimension reduction (LSDR) is a recently proposed SDR method that can overcome the limitations of KDR (Suzuki & Sugiyama, 2010). That is, LSDR is equipped with a natural model selection procedure based on cross-validation.

The proposed SCA can actually be regarded as a computationally efficient alternative to LSDR. Indeed, LSDR can also be interpreted as a dependence estimation-maximization algorithm (see Section 2.2), and the dependence estimation procedure is essentially the same as the proposed SCA, i.e., LSMI is used. The dependence maximization procedure is different from SCA—LSDR uses a *natural gradient* method (Amari, 1998).

In LSDR, the following SMI estimator is used:

$$\widetilde{\text{SMI}} = \hat{\alpha}^\top \hat{\mathbf{h}} - \frac{1}{2} \hat{\alpha}^\top \widehat{\mathbf{H}} \hat{\alpha} - \frac{1}{2},$$

where $\hat{\alpha}$, $\hat{\mathbf{h}}$ and $\widehat{\mathbf{H}}$ are defined in Section 3.1. Then the gradient of $\widetilde{\text{SMI}}$ is given by

$$\begin{aligned} \frac{\partial \widetilde{\text{SMI}}}{\partial \mathbf{W}_{\ell, \ell'}} &= \frac{\partial \hat{\mathbf{h}}^\top}{\partial \mathbf{W}_{\ell, \ell'}} (2\hat{\alpha} - \hat{\beta}) - \hat{\alpha}^\top \frac{\partial \widehat{\mathbf{H}}}{\partial \mathbf{W}_{\ell, \ell'}} \left(\frac{3}{2} \hat{\alpha} - \hat{\beta} \right) \\ &\quad + \hat{\alpha}^\top \frac{\partial \mathbf{R}}{\partial \mathbf{W}_{\ell, \ell'}} (\hat{\beta} - \hat{\alpha}), \end{aligned}$$

where $\hat{\beta} = (\widehat{\mathbf{H}} + \lambda \mathbf{R})^{-1} \widehat{\mathbf{H}} \hat{\alpha}$. The *natural gradient* update of \mathbf{W} , which takes into account the structure of the Stiefel manifold (Amari, 1998), is given by

$$\mathbf{W} \leftarrow \mathbf{W} \exp \left(\eta \left(\mathbf{W}^\top \frac{\partial \widetilde{\text{SMI}}}{\partial \mathbf{W}} - \frac{\partial \widetilde{\text{SMI}}}{\partial \mathbf{W}}^\top \mathbf{W} \right) \right),$$

where ‘exp’ for a matrix denotes the *matrix exponential*. $\eta \geq 0$ is a step size, which may be optimized by a line-search method such as *Armijo’s rule* (Patriksson, 1999).

Since cross-validation is available for model selection of LSMI, LSDR is more favorable than KDR. However, its optimization still relies on a gradient-based method and thus it is computationally expensive.

Furthermore, there seems no good initialization scheme of the transformation matrix \mathbf{W} . In the original paper by Suzuki & Sugiyama (2010), initial values were chosen randomly and the gradient method was run many times for finding a better local solution.

The proposed SCA method can successfully overcome the above weaknesses of LSDR, by providing an analytic-form solution (see Section 3.2) and a systematic initialization scheme (see Section 3.3).

5. Experiments

In this section, we experimentally investigate the performance of the proposed and existing SDR methods using artificial and real-world datasets.

5.1. Artificial Datasets

We use four artificial datasets, and compare the proposed SCA, LSDR¹ (Suzuki & Sugiyama, 2010), KDR² (Fukumizu et al., 2009), sliced inverse regression (SIR)³ (Li, 1991), sliced average variance estimation (SAVE)³ (Cook, 2000), and principal Hessian direction (pHd)³ (Li, 1992).

In SCA, we use the Gaussian kernel for \mathbf{y} :

$$L(\mathbf{y}, \mathbf{y}_\ell) = \exp \left(-\frac{\|\mathbf{y} - \mathbf{y}_\ell\|^2}{2\sigma_y} \right).$$

The identity matrix is used as regularization matrix \mathbf{R} , and the kernel widths σ_x , σ_y , and σ_z as well as the regularization parameter λ are chosen based on 5-fold cross-validation.

The performance of each method is measured by

$$\frac{1}{\sqrt{2m}} \|\widehat{\mathbf{W}}^\top \widehat{\mathbf{W}} - \mathbf{W}^{*\top} \mathbf{W}^*\|_{\text{Frobenius}},$$

where $\|\cdot\|_{\text{Frobenius}}$ denotes the Frobenius norm, $\widehat{\mathbf{W}}$ is an estimated transformation matrix, and \mathbf{W}^* is the

¹<http://sugiyama-www.cs.titech.ac.jp/~sugi/software/LSDR/index.html>

²We used the program code provided by one of the authors of Fukumizu et al. (2009), which ‘anneals’ the Gaussian kernel width over gradient iterations.

³<http://mirrors.dotsrc.org/cran/web/packages/dr/index.html>

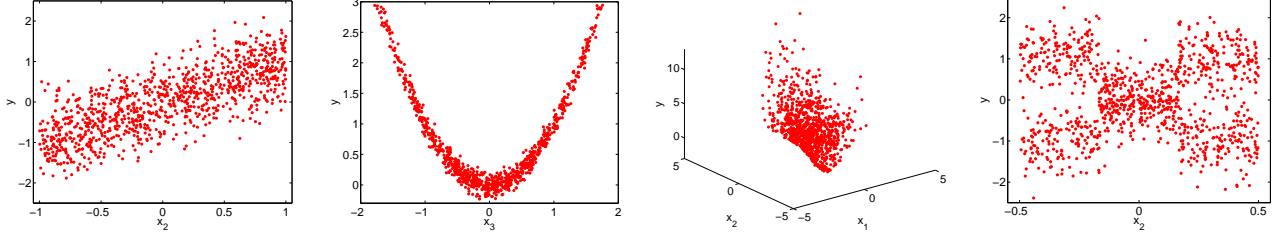


Figure 1. Artificial datasets.

Table 1. Mean of Frobenius-norm error (with standard deviations in brackets) and mean CPU time over 100 trials. Computation time is normalized so that LSDR is one. LSDR was repeated 5 times with random initialization and the transformation matrix with the minimum CV score was chosen as the final solution. ‘SCA(0)’ indicates the performance of the initial transformation matrix obtained by the method described in Section 3.3. The best method in terms of the mean Frobenius-norm and comparable methods according to the *t*-test at the significance level 1% are specified by bold face.

Datasets	<i>d</i>	<i>m</i>	SCA(0)	SCA	LSDR	KDR	SIR	SAVE	pHd
Data1	4	1	.089(.042)	.048(.031)	.056(.021)	.048(.019)	.257(.168)	.339(.218)	.593(.210)
Data2	10	1	.078(.019)	.007(.002)	.039(.023)	.024(.007)	.431(.281)	.348(.206)	.443(.222)
Data3	4	2	.065(.035)	.018(.010)	.090(.069)	.029(.119)	.362(.182)	.343(.213)	.437(.231)
Data4	5	1	.118(.046)	.042(.030)	.151(.296)	.118(.238)	.421(.268)	.356(.197)	.591(.205)
Time			0.03	0.49	1.0	0.96	<0.01	<0.01	<0.01

optimal transformation matrix. Note that the above error measure takes its value in $[0, 1]$.

We use the following four datasets (see Figure 1):

(a) **Data1:**

$$Y = X_2 + 0.5E,$$

where $(X_1, \dots, X_4)^\top \sim U([-1, 1]^4)$ and $E \sim N(0, 1)$. Here, $U(\mathcal{S})$ denotes the uniform distribution on \mathcal{S} , and $N(\mu, \Sigma)$ is the Gaussian distribution with mean μ and variance Σ .

(b) **Data2:**

$$Y = (X_3)^2 + 0.1E,$$

where $(X_1, \dots, X_{10})^\top \sim N(\mathbf{0}_{10}, \mathbf{I}_{10})$ and $E \sim N(0, 1)$.

(c) **Data3:**

$$Y = \frac{(X_1)^2 + X_2}{0.5 + (X_2 + 1.5)^2} + (1 + X_2)^2 + 0.1E,$$

where $(X_1, \dots, X_4)^\top \sim N(\mathbf{0}_4, \mathbf{I}_4)$ and $E \sim N(0, 1)$.

(d) **Data4:**

$$Y|X_2 \sim \begin{cases} N(0, 0.2) & \text{if } X_2 \leq |1/6| \\ 0.5N(1, 0.2) & \text{otherwise} \\ +0.5N(-1, 0.2), & \end{cases}$$

where $(X_1, \dots, X_5)^\top \sim U([-0.5, 0.5]^5)$ and $E \sim N(0, 1)$.

The performance of each method is summarized in Table 1, which depicts the mean and standard deviation of the Frobenius-norm error over 100 trials when the number of samples is $n = 1000$. As can be observed, the proposed SCA overall performs well. ‘SCA(0)’ in the table indicates the performance of the initial transformation matrix obtained by the method described in Section 3.3. The result shows that SCA(0) gives a reasonably good transformation matrix with a tiny computational cost. Note that KDR and LSDR have high standard deviation for Data3 and Data4, meaning that KDR and LSDR sometimes perform poorly.

5.2. Multi-label Classification for Real-world Datasets

Finally, we evaluate the performance of the proposed method in real-world multi-label classification problems.

5.2.1. SETUP

Below, we compare SCA, Multi-label Dimensionality reduction via Dependence Maximization (MDDM)⁴ (Zhang & Zhou, 2010), Canonical Correlation Anal-

⁴<http://cs.nju.edu.cn/zhouch/zhouch.files/publication/annex/MDDM.htm>

ysis (CCA)⁵ (Hotelling, 1936), and Principal Component Analysis (PCA)⁶ (Bishop, 2006). We use a real-world image classification dataset called the *PASCAL Visual Object Classes (VOC) 2010* dataset (Everingham et al., 2010) and a real-world automatic audio-tagging dataset called the *Freesound* dataset (The Freesound Project, 2011). Since the computational costs of KDR and LSDR were unbearably large, we decided not to include them in the comparison.

We employ the misclassification rate by the nearest-neighbor classifier as a performance measure:

$$\text{err} = \frac{1}{nc} \sum_{i=1}^n \sum_{k=1}^c I(\hat{y}_{i,k} \neq y_{i,k}),$$

where c is the number of classes, \hat{y} and y are the estimated and true labels, and $I(y \neq y')$ is the indicator function.

For SCA and MDDM, we use the following kernel function (Sarwar et al., 2001) for \mathbf{y} :

$$L(\mathbf{y}, \mathbf{y}') = \frac{(\mathbf{y} - \bar{\mathbf{y}})^\top (\mathbf{y}' - \bar{\mathbf{y}})}{\|\mathbf{y} - \bar{\mathbf{y}}\| \|\mathbf{y}' - \bar{\mathbf{y}}\|},$$

where $\bar{\mathbf{y}}$ is the sample mean: $\bar{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i$.

5.2.2. PASCAL VOC 2010 DATASET

The VOC 2010 dataset consists of 20 binary classification tasks of identifying the existence of a person, aeroplane, etc. in each image. The total number of images in the dataset is 11319, and we used 1000 randomly chosen images for training and the rest for testing.

In this experiment, we first extracted visual features from each image using the *Speed Up Robust Features* (SURF) algorithm (Bay et al., 2008), and obtained 500 *visual words* as the cluster centers in the SURF space. Then, we computed a 500-dimensional *bag-of-feature* vector by counting the number of visual words in each image. We randomly sampled the training and test data 100 times, and computed the means and standard deviations of the classification error.

The results are plotted in Figure 2(a), showing that SCA outperforms the existing methods, and SCA is the only method that outperforms ‘ORI’ (no dimension reduction)—SCA achieves almost the same error rate as ‘ORI’ with only a 10-dimensional subspace.

5.2.3. FREESOUND DATASET

The *Freesound* dataset (The Freesound Project, 2011) consists of various audio files annotated with word tags

such as ‘people’, ‘noisy’, and ‘restaurant’. We used 230 tags in this experiment. The total number of audio files in the dataset is 5905, and we used 1000 randomly chosen audio files for training and the rest for testing.

We first extracted *Mel-Frequency Cepstrum Coefficients* (MFCC) (Rabiner & Juang, 1993) from each audio file, and obtained 1024 *audio features* as the cluster centers in MFCC. Then, we computed a 1024-dimensional *bag-of-feature* vector by counting the number of audio features in each audio file. We randomly chose the training and test samples 100 times, and computed the means and standard deviations of the classification error.

The results plotted in Figure 2(b) show that, similarly to the image classification task, the proposed SCA outperforms the existing methods, and SCA is the only method that outperforms ‘ORI’.

6. Conclusion

In this paper, we proposed a novel *sufficient dimension reduction* (SDR) method called *sufficient component analysis* (SCA), which is computationally more efficient than existing SDR methods. In SCA, a transformation matrix was estimated by iteratively performing dependence estimation and maximization, both of which are *analytically* carried out. Moreover, we developed a systematic method to design a good initial transformation matrix, which highly contributes to further reducing the computational cost and help obtain a good local optimum solution. We applied the proposed SCA to real-world image classification and audio tagging tasks, and experimentally showed that the proposed method is promising.

Acknowledgments

The authors thank Prof. Kenji Fukumizu for providing us the KDR code and Prof. Taiji Suzuki for his valuable comments. MY was supported by the JST PRESTO program. GN was supported by the MEXT scholarship. MS was supported by SCAT, AOARD, and the JST PRESTO program.

References

- Amari, S. Natural gradient works efficiently in learning. *Neural Computation*, 10:251–276, 1998.
- Bay, H., Ess, A., Tuytelaars, T., and Gool, L. V. Surf: Speeded up robust features. *Computer Vision and Image Understanding*, 110(3):346–359, 2008.
- Bishop, C. M. *Pattern Recognition and Machine*

⁵<http://www.mathworks.com/help/toolbox/stats/canoncorr.html>

⁶<http://www.mathworks.com/help/toolbox/stats/princomp.html>

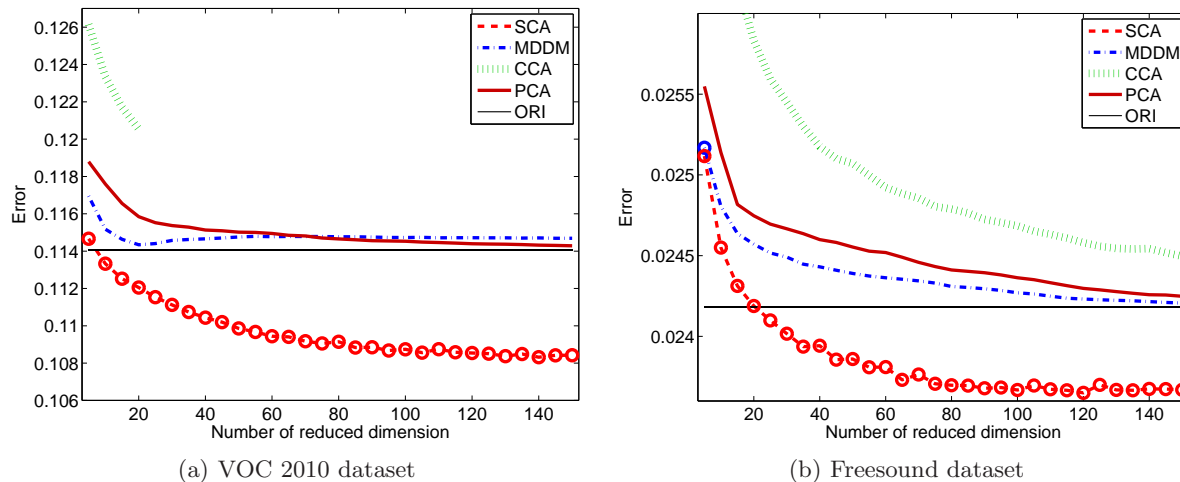


Figure 2. Results on image classification with VOC 2010 dataset and audio classification with Freesound datasets. Misclassification rate when the one-nearest-neighbor classifier is used as a classifier is reported. The best dimension reduction method in terms of the mean error and comparable methods according to the t-test at the significance level 1% are specified by ‘o’. CCA can be applied to dimension reduction up to c dimensions, where c is the number of classes ($c = 20$ in VOC 2010 and $c = 230$ in Freesound). ‘ORI’ denotes the original data without dimension reduction.

Learning. Springer, New York, NY, 2006.

Cook, R. D. *Regression graphics: Ideas for studying regressions through graphics*. Wiley, New York, 1998.

Cook, R. D. Save: A method for dimension reduction and graphics in regression. *Theory and Methods*, 29: 2109–2121, 2000.

Epanechnikov, V. Nonparametric estimates of a multivariate probability density. *Theory of Probability and its Applications*, 14:153–158, 1969.

Everingham, M., Gool, L. V., Williams, C. K. I., Winn, J., and Zisserman, A. The PASCAL Visual Object Classes Challenge 2010 (VOC2010) Results. <http://www.pascal-network.org/challenges/VOC/voc2010/workshop/index.html>, 2010.

Fukumizu, K., Bach, F. R., and Jordan, M. Kernel dimension reduction in regression. *The Annals of Statistics*, 37(4):1871–1905, 2009.

Hotelling, H. Relations between two sets of variates. *Biometrika*, 28:321–377, 1936.

Li, K.-C. Sliced inverse regression for dimension reduction. *Journal of American Statistical Association*, 86:316–342, 1991.

Li, K.-C. On principal Hessian directions for data visualization and dimension reduction: Another application of Stein’s lemma. *Journal of American Statistical Association*, 87:1025–1034, 1992.

Patriksson, M. *Nonlinear Programming and Variational Inequality Problems*. Kluwer Academic, Dordrecht, 1999.

Rabiner, L. and Juang, B.-H. *Fundamentals of Speech Recognition*. Prentice Hall, Englewood Cliffs, NJ, 1993.

Sarwar, B., Karypis, G., Konstan, J., and Reidl, J. Item-based collaborative filtering recommendation algorithms. In *Proceedings of the 10th international conference on World Wide Web (WWW2001)*, pp. 285–295, 2001.

Suzuki, T. and Sugiyama, M. Sufficient dimension reduction via squared-loss mutual information estimation. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS2010)*, pp. 804–811, 2010.

Suzuki, T., Sugiyama, M., Kanamori, T., and Sese, J. Mutual information estimation reveals global associations between stimuli and biological processes. *BMC Bioinformatics*, 10(S52), 2009.

The Freesound Project. Freesound, 2011. <http://www.freesound.org>.

Zhang, Y. and Zhou, Z.-H. Multilabel dimensionality reduction via dependence maximization. *ACM Trans. Knowl. Discov. Data*, 4:14:1–14:21, 2010. ISSN 1556-4681.